

Assignment 2: Variational Inference in the TrueSkill model

- Name:
- StudentNumber:
- Collaborators:

Goal

The goal of this assignment is to become familiar with the basics of Bayesian inference in large models with continuous latent variables, and implement stochastic variational inference with Automatic Differentiation.

Outline

1. Introduce our problem and describe our model.
2. Implement our model.
3. Investigate the model on easy-to-visualize toy data.
4. Implement Stochastic Variational Inference (SVI) objective (ELBO) manually, using Automatic Differentiation (e.g. **Zygote.jl**).
5. Use our SVI implementation to learn an approximate distribution inferred from the toy data and visualize.
6. Use SVI to perform approximate inference on real data. A collection of games played by by Woman Grandmasters on chess.com. The data was modified from **this Kaggle Dataset**.
7. Use variational approximation to estimate inference questions using the model.

Background

We will implement a variant of the TrueSkill model, a player ranking system for competitive games originally developed for Halo 2. It is a generalization of the Elo rating system in Chess.

Here are some readings you can familiarize yourself with the problem and their model. Note that we will use different inference methods than those discussed in these works (message passing):

- **The 2006 technical report at Microsoft Research**
- **The 2007 NeurIPS conference paper introducing TrueSkill**
- **The 2018 followup work TrueSkill 2**

Our assignment is based on **one developed by Carl Rasmussen at Cambridge for his course on probabilistic machine learning**. In their assignment they implement and utilize a Gibbs sampling method for approximate inference. We will not implement an MCMC sampling method, instead using SVI for our approximate inference.

1. Problem and Model Description

We will consider a simplified version of the TrueSkill model to model the skill of a individual players i at a 2-player game: chess.

We assume that each player has an unknown skill $z_i \in \mathbb{R}$. These are called latent variables to our model since we cannot observe the players' skill directly.

Given our list of players I , we have no way *a priori* (before we observe data) to infer the players' unknown skill. If we were familiar with the roster of Woman Grandmaster (WGM) chess players, though, we could incorporate our *prior information* about the players to suggest plausible skills. However, these initial priors can be biased, e.g. player fame, charisma, playstyle.

Instead, we will simply assume that all players' skill is distributed according to a standard Gaussian distribution.

We also assume that all players' skills are *a priori* independent.

Because we cannot monitor players' skill directly, we must collect other evidence. We observe the players' performance against other players in a series of games. We will use the observed game data to update our model of the players' skill.

Our model of the likelihood that player i beats player j in a game of chess depends on our model of their respective skills z_i, z_j . We will use the following likelihood model:

$$p(\text{game}(i \text{ wins}, j \text{ loses}) \mid z_i, z_j) = \sigma(z_i - z_j)$$

where

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

Note that $\log(1 + \exp(y))$ can suffer significant numerical instability if not implemented carefully. You should use a function like `StatsFuns.log1pexp` or its `numpy` equivalent.

We also assume that each game is independent from the others, given the players' skill.

Data

The observations of game outcomes are collected in an array. Each column contains a pair of player indices, i, j for the chess game. The first index, player i is the winner, and so j lost the game.

If we observe M games played then the data is collected into an array with shape $2 \times M$.

Note that data retains its dimension even as we observe more players. Additional players are identified by their index. If we wanted to consider multiplayer games, like Halo 2, we could increase the first dimension, from our 2-player chess game.

- using `StatsFuns.log1pexp`

2. Implementing the Model

- log-prior of skills z_s
- log-likelihood of observing game with player i beating player j
- log-likelihood of observing collection of games with players' skills z_s
- log-likelihood of joint distribution observed games and skills z_s

```
• function log_prior(zs)
•   return #TODO
• end
```

```
• function logp_i_beats_j(zi,zj)
•   return #TODO
• end
```

```
• function all_games_log_likelihood(zs,games)
•   zs_a = #TODO
•   zs_b = #TODO
•   likelihoods = #TODO
•   return #TODO
• end
```

```
• function joint_log_density(zs,games)
•   #hint: you have likelihood and prior, use Bayes
•   return #TODO
• end
```

3. Visualize the Model on Toy Data

To check our understanding of the data and our model we consider a simple scenario.

Toy Data

Let's model the chess skills for *only* two players, A and B . Restricting to 2-dimensions allows us to visualize the *posterior* distribution of our model. Unlike the high-dimensional distributions we would model with real data, we can evaluate the model on a grid of points, like numerical integration, to produce a plot of contours for the possibly complicated posterior.

We provide a function `two_player_toy_games` which produces toy chess data for two players. e.g. `two_player_toy_games(5,3)` produces a dataset where player A wins 5 games and player B wins 3 games.

2D Posterior Visualization

You can use the function `skillcontour!` to perform this grid of evaluations over the 2D latent space $z_A \times z_B$. This will plot isocontours, curves of equal density, for the posterior distribution.

As well, `plot_line_equal_skill!` is provided to simply indicate the region of latent space corresponding to the players having equal skill, $z_A = z_B$.

(the default settings for these plotting functions may need modification.)

- `two_player_toy_games(p1_wins, p2_wins) = cat(collect.([repeat([1,2]',p1_wins)', repeat([2,1]',p2_wins)'])), dims=2)`

- `two_player_toy_games(5,3)`

- `using Plots`

```
function skillcontour!(f; colour=nothing)
    n = 100
    x = range(-3,stop=3,length=n)
    y = range(-3,stop=3,length=n)
    z_grid = Iterators.product(x,y) # meshgrid for contour
    z_grid = reshape.(collect.(z_grid),:,1) # add single batch dim
    z = f.(z_grid)
    z = getindex.(z,1)
    max_z = maximum(z)
    levels = [.99, 0.9, 0.8, 0.7,0.6,0.5, 0.4, 0.3, 0.2] .* max_z
    if colour==nothing
        p1 = contour!(x, y, z, fill=false, levels=levels)
    else
        p1 = contour!(x, y, z, fill=false, c=colour, levels=levels,colorbar=false)
    end
    plot!(p1)
end
```

```

• function plot_line_equal_skill!()
•   plot!(range(-3, 3, length=200), range(-3, 3, length=200), label="Equal Skill")
• end

```

Plot the following distributions for two players A and B . In each, also plot the line of equal skill $z_A = z_B$.

Note, before plotting convert from log-density to density. Since there is only 1 distribution per plot, we are not comparing distributions, we do not need to worry about normalization of the contours.

1. Isocontours of the prior distribution over players' skills.
2. Isocontours of the likelihood function.
3. Isocontours of the posterior over players' skills given the observaton: player A beat player B in 1 game .
4. Isocontours of the posterior over players' skills given the observation: player A beat player B in 10 games .
5. Isocontours of the posterior over players' skills given the observation: 20 games were played, player A beat player B in 10 games .

```

• #1
• # TODO: plot prior contours

```

```

• #2
• # TODO: plot likelihood contours

```

```

• #3
• # TODO: plot posterior contours with player A winning 1 game

```

```

• #4
• # TODO: plot joint contours with player A winning 10 games

```

```

• #5
• #TODO: plot joint contours with player A winning 10 games and player B winning 10 games

```

4. Stochastic Variational Inference with Automatic Differentiation

A nice quality of our Bayesian approach is that we separate the model specification from the (approximate) method we use for inference.

In the original TrueSkill paper they described a message passing strategy for approximate inference. In Carl Rasmussen's assignment the students implement Gibbs sampling, a kind of Markov Chain Monte Carlo method.

We will use gradient-based stochastic variational inference (SVI), a recent technique that is extremely successful in our domain.

We will use SVI to produce an approximate posterior distribution to the, possibly complicated, posterior distributions specified by our model.

1. Implement the function `elbo` which computes an unbiased estimate of the Evidence Lower Bound. The ELBO is proportional to a KL divergence between the model posterior $p(z \mid \text{data})$ and the approximate variational distribution $q_\phi(z \mid \text{data})$.
2. Implement an objective function to optimize the parameters of the variational distribution. We will minimize a "loss" with gradient descent, the negative ELBO estimated with 100 samples.
3. Write an optimization procedure `learn_and_vis_toy_variational_approx` that takes initial variational parameters and observed data and performs a number of iterations of gradient optimization where each iteration:
 1. compute the gradient of the loss with respect to the variational parameters using AD
 2. update the variational parameters taking a `lr`-scaled step in the direction of the descending gradient.
 3. report the loss with the new parameters (using `@info` or `print` statements)
 4. on one plot axis plot the target posterior contours in red and the variational approximation contours in blue. `display` the plot or save to a folder.

```
function elbo(params, logp, num_samples)
    samples = #TODO
    logp_estimate = #TODO
    logq_estimate = #TODO
    return #TODO: should return scalar (hint: average over batch)
end
```

```
# Convenience function for taking gradients
function neg_elbo(params; games = two_player_toy_games(1,0), num_samples = 100)
    # TODO: Write a function that takes parameters for q,
    # evidence as an array of game outcomes,
    # and returns the -elbo estimate with num_samples many samples from q
    logp(zs) = joint_log_density(zs, games)
    return -elbo(params, logp, num_samples)
end
```

```
function learn_and_vis_toy_variational_approx(init_params, toy_evidence; num_itrs=200,
lr= 1e-2, num_q_samples = 10)
    params_cur = init_params
    for i in 1:num_itrs
        grad_params = #TODO: gradients of variational objective with respect to parameters
        params_cur = #TODO: update parameters with lr-sized step in descending gradient
        #TODO: report the current elbbo during training
        # TODO: plot true posterior in red and variational in blue
        # hint: call 'display' on final plot to make it display during training
        plot();
        #TODO: skillcontour!(..., colour=:red) plot likelihood contours for target
        posterior
        # plot_line_equal_skill()
```

- *#TODO: display(skillcontour!(..., colour=:blue)) plot likelihood contours for variational posterior*
- end
- return params_cur
- end

5. Visualizing SVI on Two Player Toy

For the following collections of observed toy data of two player games learn an variational distribution to approximate the model posterior.

Using the plots produced in the training SVI procedure, one plot per iteration, create an animation that shows the variational distribution converging to its final approximation. In your final report, please simply plot the the learned variational distribution in at the final iteration.

Each of these require its own set of initial variational parameters and variational optimization procedure. The plots should have contours for two distributions, model posterior (red) and variational approximation (blue). To allow useful visualization do **not** normalize the contours to the same scale.

1. Report the final loss and plot the posteriors for the observation: player A beats player B in 1 game
2. Report the final loss and plot the posteriors for the observation: player A beats player B in 10 games
3. Report the final loss and plot the posteriors for the observation: player A wins 10 games player B wins 10 games

toy_mu =

- *# Toy game*
- toy_mu = [-2.,3.] *# Initial mu, can initialize randomly!*

toy_ls =

- toy_ls = [0.5,0.] *# Initial log-sigma, can initialize randomly!*

toy_params_init =

- toy_params_init = (toy_mu, toy_ls)

- *#1*
- *#TODO: fit q with SVI observing player A winning 1 game*
- *#TODO: save final posterior plots*

- *#2*
- *#TODO: fit q with SVI observing player A winning 10 games*
- *#TODO: save final posterior plots*

- *#3*

- *#TODO: fit q with SVI observing player A winning 10 games and player B winning 10 games*
- *#TODO: save final posterior plots*

6. Approximate Inference on Real Data

Original Dataset

We will model the skills of players on chess.com from a collection of games played by Women GrandMasters since September 2009.

I have modified the dataset collected from the Kaggle Dataset: **Chess Games of Woman Grandmasters (2009 - 2021)**.

Note the extra information contained in the original dataset of game observations. If we were familiar with the rules of chess when we designed our model, and aware of this extra data collected in our observations, we could describe a better model. For example,

- since White moves first, the White player has a considerable advantage
- not all games are observed as win, lose, some games end in draw due to repetition, or win by timeout.
- games have different timing rules.

Our model is simple and does not make use of that extra information in the original dataset. For this assignment it will not be necessary to modify our model beyond the simple case. If you are interested in extending this, consider how you could adapt the model specification to incorporate these observations.

Modified Dataset

My modifications to the data remove extra information about which player is White and simplifies the win conditions.

The data consists of two arrays:

- `games` is a $2 \times M$ collection of chess game outcomes, one game per column. The first row contains the indices of players who won, the second contains the indices of players who lost.
- `names` is a vector of player names, `names[i]` corresponds to the player chess.com username indicated by the values in `games`.

- using CSV, DataFrames

`games =`


```
• games = collect(Array(CSV.read("games.csv", DataFrame)))'
```

```
names =
```

```
• names = vec(Array(CSV.read("names.csv", DataFrame)))
```

Variational Distribution

Use a fully-factorized Gaussian distribution for $q_\phi(z \mid data)$.

$$q_\phi(z \mid games) = \prod_i \mathcal{N}(z_i \mid \mu_i, \sigma_i)$$

Note, for numerical stability we work with the `log_density`. Also, for unconstrained optimization, we parameterize our Gaussian with $\log \sigma$, not variance or standard deviation.

You will need to implement this variational distribution with parameters $\phi = [\mu, \log \sigma]$ by defining `logq(z, params)` inside the `elbo`.

Using the model and SVI method, we will condition our posterior on observations from the dataset.

In the previous Two Player Toy we were informed about the players' skill by observing games between them. Now we have many players playing games among eachother. For any two players, i and j , answer yes or no to the following:

1. In general, is $p(z_i, z_j \mid \text{all games})$ proportional to $p(z_i, z_j, \text{all games})$?
2. In general, is $p(z_i, z_j \mid \text{all games})$ proportional to $p(z_i, z_j \mid \text{games between } i \text{ and } j)$? That is, do the games between player i and j provide all the information about the skills z_i and z_j ?

Hint: consider the graphical model for three players, i, j, k and condition on results of games between all 3. Examine the conditional independences in the model. (graphical model is not required for marks)

3. write a new optimization procedure `learn_variational_approx` like the one from the previous question, but **does not produce any plots**. As well, have the optimization procedure save the loss (negative ELBO) for each iteration into an array. Initialize a variational distribution and optimize it to learn the model joint distribution observing the chess games in the dataset.
4. Plot the losses (negative ELBO estimate) obtained during variational optimization and report the loss of the final approximation.
5. We now have an approximate posterior over all our players' skill. Our variational family is simple, allowing us to easily compute an approximate mean and variance for all players' skill. Sort the players by mean skill and list the names of the 10 players with highest mean skill under our variational approximation. Use `sortperm`.
6. Plot the mean and variance of all players' skill, sorted by mean skill. There is no need to include the names of the players in the x-axis. Use `plot(means, yerror=exp.(logstds))` to add

variance to the plot.

```

• function learn_variational_approx(init_params, tennis_games; num_itrs=200, lr= 1e-2,
  num_q_samples = 10)
•   params_cur = init_params
•   for i in 1:num_itrs
•     grad_params = #TODO: gradients of variational objective wrt params
•     params_cur = #TODO: update parameters with lr-sized steps in descending gradient direction
•     log_loss #TODO: save objective value with current parameters
•   end
•   return params_cur
• end

```

- *#1 Yes or No?*
- *#2 Yes or No?*

`init_params =`

- *#3*
- *# TODO: Initialize variational family*
- `init_params = (init_mu, init_log_sigma)` *# random initialization*

`variational_params =`

- `variational_params = learn_variational_approx(init_params, games)`

- *#4*
- *#TODO plot losses during ELBO optimization and report final loss*

- *#5*
- *#TODO: sort players by mean skill under our model and list top 10 players.*

- *#6*
- *#TODO: plot mean and variance of all players, sorted by skill*

7. More Approximate Inference with our Model

Some facts from [The Queens of Chess Kaggle Notebook](#) motivate a couple inference questions.

- `camillab` Camilla Baginskaite is a Lithuanian and American WGM and is the most active player in the dataset, over 60K games played in the last 12 years!
- `liverpoolborn` Sheila Jackson is an English WGM and the most successful player (highest win %).
- `meri_arabidze` Meri Arabidze has achieved the highest rating by a WGM of 2763 on chess.com.
- `sylvanaswindrunner` Jovana Rapport is currently rated 2712, very close to Meri's all time high!

Here I find these players in the `names` data to determine their indices.

`camillab =`

```
• camillab = findfirst(i -> i == "camillab", names)
```

```
• length(findall(g -> g == 4059, vec(games))) # number of games camillab played!
```

```
liverpoolborn =
```

```
• liverpoolborn = findfirst(i -> i == "liverpoolborn", names)
```

```
meri_arabidze =
```

```
• meri_arabidze = findfirst(i -> i == "meri-arabidze", names)
```

Use the optimized variational approximation to give estimates for the following inference problem:

What is the probability player A is better than player B?

Recall from our early plots, this is the mass under the distribution in the region above the line of equal skill, $z_A > z_B$.

We have chosen a very simple variational family, factorized Gaussians. This would allow us to derive the exact expression for the probability under the variational approximation that player A has higher skill than player B. This year we will not derive this value exactly.

We can estimate the answer to our inference question by sampling $N = 10000$ samples from the variational distribution and using a Simple Monte Carlo estimate.

This inference question restricts us back to a two-player scenario. We know from previous questions that we can plot contours for distributions over 2 player skills.

most games vs most successful

1. Plot the **approximate** posterior over the skills of camillab and liverpoolborn.
2. Estimate the probability that liverpoolborn is more skillful than camillab.

all time high vs contender

3. Plot the **approximate** posterior over the skills of meri_arabidze and sylvanaswindrunner.
4. Estimate the probability that sylvanaswindrunner is more skillful than meri_arabidze.

```
• #1
• #TODO: plot approximate posterior over two players
```

```
• #2
• # TODO: use Simple Monte Carlo to estimate probability player A is better
```

```
• #3
• #TODO: plot approximate posterior over two players
```

- #4
- # *TODO: use Simple Monte Carlo to estimate probability player A is better*